

# Code Generation: Results Achieved with Automated Chain-of-Thought prompting

	DSC 1.3B	DSC 6.7B	DSC 33B	CL 7B	CL 13B	CL 34B	GPT 3.5	GPT 4	Human Written	Own <i>vs</i> LLMs	Own <i>vs</i> LLMs \ F	Own <i>vs</i> Human
DSC 1.3B	0.73	0.66	0.60	0.65	0.51	0.54	0.57	0.59	-0.25	** (N)	** (N)	*** (L)
DSC 6.7B	0.53	0.60	0.62	0.59	0.42	0.49	0.50	0.60	-0.30	(N)	(N)	*** (L)
DSC 33B	0.38	0.48	0.54	0.36	0.35	0.32	0.35	0.54	-0.38	* (N)	* (N)	*** (L)
CL 7B	0.72	0.65	0.62	0.70	0.52	0.61	0.58	0.67	-0.13	(N)	(N)	*** (L)
CL 13B	0.55	0.50	0.58	0.49	0.36	0.44	0.47	0.45	-0.19	(N)	(N)	*** (M)
CL 34B	0.43	0.40	0.42	0.46	0.33	0.44	0.37	0.41	-0.34	* (N)	* (N)	*** (L)
GPT-3.5	0.15	0.17	0.23	0.14	-0.01	0.09	0.15	0.27	-0.69	(N)	(N)	*** (L)
GPT-4	0.09	0.18	0.18	0.13	0.15	0.15	0.24	0.43	-0.68	*** (S)	*** (S)	*** (L)
Average (all)	0.32	0.35	0.39	0.31	0.23	0.29	0.32	0.42	-0.46	-	-	-
Average (large)	0.45	0.46	0.47	0.44	0.33	0.39	0.40	0.49	-0.37	-	-	-

Adjusted  $p$ -values: \* <0.05, \*\* <0.01, \*\*\* <0.001. Cliff delta: N=Negligible, S=Small, M=Medium, L=Large

Table 1: (Python) Average of differences between the LLM judgments (0 or 1) and the ground truth (*i.e.*, 1 if the method passes the test and 0 otherwise). Last three columns report adj.  $p$ -value and effect size when comparing the judgements each LLM gave to functions it generated against those it gave when judging functions (i) generated by all other LLMs, (ii) generated by all other LLMs but those belonging to the same family, and (iii) written by humans.